



JP2001067187

Biblio Page 1

esp@cenet

**STORAGE SUB-SYSTEM AND ITS CONTROL METHOD**

Patent Number: JP2001067187
Publication date: 2001-03-16
Inventor(s): ARAKAWA TAKASHI; MOGI KAZUHIKO; YAMAKAMI KENJI; ARAI HIROHARU
Applicant(s):: HITACHI LTD
Requested Patent: ☐ JP2001067187 (JP01067187)
Application Number: JP19990242713 19990830
Priority Number(s):
IPC Classification: G06F3/06 ; G06F12/00
EC Classification:
Equivalents:

Abstract

PROBLEM TO BE SOLVED: To simplify a work for optimizing arrangement by re-arrangement by the user of a disk array system or the like by changing the correspondence of a logical storage area from a physical storage area into the second physical storage area and executing re-arrangement.

SOLUTION: A control part 300 automatically executes re-arrangement execution processing at the set time and date. That is, the part 300 copies contents stored in a re-arrangement source physical area in a re-arrangement destination physical area based on re-arrangement information 408. Moreover, at the point of time when the copying is completed and the whole contents of the re-arrangement source physical area are reflected in the re-arrangement destination physical area, the control part 300 changes a physical area corresponding to a logical area for executing re-arrangement in logical/physical correspondence information 400 from the re-arrangement source physical area into the re-arrangement destination physical area. Besides, the control part 300 uses the re-arrangement destination physical area on a non-usage physical area 1470, changes the re-arrangement source physical area into the non-usage one and, moreover, updates the time and date of re-arrangement execution time information 406 into the one for a next time by referring to time and date updating information on re-arrangement execution time information 406.

Data supplied from the esp@cenet database - I2

16/10/01 x 2500

るストレージサブシステム。

【請求項10】請求項6、7、8、または9に記載のストレージサブシステムであって、ストレージサブシステムは、複数のディスク装置を有するディスクアレイであり、前記ディスク装置の使用率を使用状況情報として用いる手段を有することを特徴とするストレージサブシステム。

【発明の詳細な説明】

【0001】 発明の属する技術分野】 本発明は、複数の記憶装置を有するストレージサブシステム、およびその制御方法に関する。

【0002】

【従来の技術】 コンピュータシステムにおいて、高性能を実現する二次記憶システムの一つにディスクアレイシステムがある。ディスクアレイシステムは、複数のディスク装置をアレイ状に配置し、前記各ディスク装置に分割格納されるデータのリード/ライトを、前記各ディスク装置を並列に動作させることによって、高速に行うシステムである。ディスクアレイシステムに関する論文として、D. A. Patterson, G. Gibbs, and R. H. Kats, "A Case for Redundant Arrays of Inexpensive Disks (RAID)" (in Proc. ACM SIGMOD, pp. 109-116, June 1988) がある。この論文では、冗長性を付加したディスクアレイシステムに

対し、その構成に応じてレベル1からレベル5の区別をみている。これらの区別に加えて、冗長性無し（レイスクアレイシステム）をレベル0と呼ぶことも多い。上記の各レベルは冗長性などにより実現するためのコストや性能特性などが異なるため、ディスクアレイシステムを構築するにあたって、複数のレベルのアレイ（ディスク装置の組）を混在させることも多い。ここでは、この組

のことをパーティグループと呼ぶ。

【0003】 ディスク装置は、性能や容量などによりコストが異なり、ディスクアレイシステムを構築するにあたって最適なコストパフォーマンスを実現するために、やはり性能や容量の異なる複数のディスク装置を用いることがある。

【0004】 ディスクアレイシステムに格納されるデータを上記のようにディスク装置に分散して配置するため、ディスクアレイシステムは、ディスクアレイシステムに接続するホストコンピュータがアクセスする論理記憶領域とディスク装置の記憶領域とを対応づけて物理記憶領域への対応づけを行う。特開9-27454号公報には、ホストコンピュータからの論理記憶領域に対する1/Oアクセスについての情報を取得する手段と、論理記憶領域の物理記憶領域への対応づけを変更して物理的再配置を行う手段により、格納されたデータ

の最適配置を実現するディスクアレイシステムが開示されている。

【0005】

【発明が解決しようとする課題】 特開9-27454号公報に示されるような従来の技術における最適配置の実行方法については以下の課題がある。

【0006】 再配置する論理記憶領域の選択および再配置先の物理記憶領域の選択にあたり、ディスクアレイシステムのユーザまたは保守員が、前記ディスクアレイシステムの構成や個々のディスク装置の特性や性能などの情報を参照して前記選択を行わなければならない。ユーザまたは保守員による作業が煩雑となっていた。

【0007】 また、ディスクアレイシステムが選択を自動的に行う場合においても、ユーザまたは保守員が前記個々のディスク装置の情報を参照して選択基準値を設定しなければならない。やがてユーザまたは保守員による作業が煩雑となっていた。特に、上記のように個々のレベルや個々のディスク装置の混在するディスクアレイシステムについては情報管理の煩雑さが増大する。

【0008】 また、ディスクアレイシステムが選択のために1/Oアクセス情報の参照は、ホストコンピュータおよびディスクアレイシステムを含むシステムで行われる処理のスケジューリングの特性を考慮している場合、一般的にコンピュタシステムで行われる処理と処理に伴う1/Oは、ユーザによって作成されたスケジューリングに則って行われており、また処理および1/Oの傾向は日毎、月毎、年毎などの周期性を示す場合も多く、一般的にユーザは特定期間の処理および1/Oに関心があると考えられる。

【0009】 また上記従来技術において、再配置による性能チューニング方法については以下の課題がある。物理的再配置による性能チューニング方法は、ディスク装置、すなわち、物理記憶領域の使用状況に要する加えるものであるが、従来の技術においては、ホストコンピュータからの論理記憶領域に対する1/Oアクセスについての情報を参照するため、再配置する論理記憶領域の選択および再配置先の物理記憶領域の選択にあたり、正しい選択が行えない可能性があった。

【0010】 また、ホストコンピュータからのシーケンシャルアクセスとランダムアクセスが混在する同一のディスク装置に含まれる別々の物理記憶領域に対して行われる場合でも、シーケンシャルアクセスとランダムアクセスを異なるディスク装置に分離するために、再配置先の物理記憶領域を任意に特定して自動的に再配置を行わせることはできなかった。一般に、ホストコンピュータからの処理要件として、データ長の小さいランダムアクセスには短時間で応答（高応答性）が求められるが、同一ディスク装置にデータ長の大きいシーケンシャルアクセスが存在する場合、ランダムアクセスの応答時間はシーケンシャルアクセスの処理に阻害されて長くなり、

応答性能は悪化してしまう。

【0011】 本発明の第一の目的は、ディスクアレイシステムにユーザまたは保守員が再配置による最適配置を行うための作業を簡便にすることにある。

【0012】 本発明の第二の目的は、ホストコンピュータおよびディスクアレイシステムを含むシステムでの処理のスケジューリングを考慮した再配置による最適配置を実現することにある。

【0013】 本発明の第三の目的は、再配置する論理記憶領域の選択および再配置先の物理記憶領域の選択にあたり、実際の記憶装置であるディスク装置の使用状況に基づいて選択を行う。ディスクアレイシステムの制御方法およびディスクアレイシステムを提供することにある。

【0014】 本発明の第四の目的は、ディスクアレイシステムにおける同一ディスク装置での異なるシーケンシャルアクセスとランダムアクセスの混在に対し、再配置先のディスク装置を任意に特定して再配置によりシーケンシャルアクセスおよびランダムアクセスを異なるディスク装置に自動的に分離することができようすることにある。

【0015】 問題を解決するための手段】 上記の第一の目的を実現するために、1台以上のホストコンピュータに接続するディスクアレイシステムは、配下の複数のディスク装置の使用状況情報を取得する手段と、ホストコンピュータがリード/ライト対象とする論理記憶領域とディスク装置の第一の物理記憶領域との対応づけを行う手段とを有し、さらに、複数のディスク装置をそれぞれ独立性を有する複数の組（クラス）として管理する手段と、使用状況情報およびクラス属性に基づき論理記憶領域に好適な再配置先のクラスを決定する手段と、論理記憶領域の再配置先として利用可能な第二の物理記憶領域をクラス内から選択する手段と、第一の物理記憶領域の内容を前記第二の物理記憶領域にコピーするとともに論理記憶領域の対応づけを第一の物理記憶領域から第二の物理記憶領域へ変更して再配置を行う手段を備える。

【0016】 また、上記第二の目的を実現するために、ディスクアレイシステムは、使用状況情報を蓄積し、設定された期間の使用状況情報に基づき、論理記憶領域の再配置先を決定する手段と、設定された時間に再配置を行う手段を備えることができる。

【0017】 また、上記第三の目的を実現するために、ディスクアレイシステムは、使用状況情報として、ディスク装置の単位時間当たりの使用時間（使用率）を用いる手段を備える。

【0018】 また、上記第四の目的を実現するために、ディスクアレイシステムは、各クラスに属性として設定された対象アクセス種別（シーケンシャル/ランダム/アクセス種別）と使用率上限値を用いて、クラスの使用率上限値を超えている記憶装置から再配置する論理記憶領域

域を選択し、論理記憶領域に対するアクセス種別の分析結果に基づいて論理記憶領域の再配置先のクラスを好適なアクセス種別のクラスから、各クラスの使用率上限値を超えないように決定する手段を備える。

【0019】

【発明の実施の形態】 以下、本発明の実施の形態を図1〜図2を用いて説明する。

【0020】 第一の実施の形態は、再配置の形態では、クラス6000に基づく再配置の判断と、再配置判断および実行のスケジュールリングについて説明する。

【0021】 図1は、本発明の第一の実施の形態における計算機システムの構成図である。

【0022】 本実施の形態における計算機システムは、ホスト100、ストレージサブシステム200、制御部300、およびディスクアレイシステム400を有している。

【0023】 ホスト100は、ストレージサブシステム200に1/Oバス800を介して接続し、ストレージサブシステム200に対してリード/ライトの1/Oを行う。1/Oの際、ホスト100は、ストレージサブシステム200の記憶領域について論理領域を指定する。1/Oバス800の例としては、ESCON、SCSI、ファイバチャネルなどがある。

【0024】 ストレージサブシステム200は、制御部300および複数の記憶装置500を有する。制御部300は、リード/ライト処理310、使用状況情報取得処理311、再配置判断処理312、及び再配置実行処理313を行う。また、ストレージサブシステム200は、論理/物理対応情報400、クラス属性情報401、クラス属性情報402、論理領域使用状況情報403、物理領域使用状況情報404、再配置判断対象期間情報405、再配置実行時刻情報406、未使用領域情報407、及び再配置情報408を保持する。

【0025】 ホスト100、制御部300、および制御部300は、ネットワーク900で接続される。ネットワーク900の例としては、FDDI、ファイバチャネルなどがある。

【0026】 ホスト100、制御部300、および制御部300は、各々での処理を行うためのメモリ、CPUなど、計算機において一般に用いられる構成要素もそれぞれ存在するが、本実施の形態の説明においては重要であるため、ここでは説明を省略する。

【0027】 ホスト100が、ストレージサブシステム200に対してリード/ライトを行う場合のリード/ライト処理310、および使用状況情報取得処理311について、図2で説明する。

【0028】 リード/ライト処理310において、ホスト100は、ストレージサブシステム200の制御部300に対してリードまたはライトを論理領域を指定して要求する（ステップ1000）。要求を受領した制御部300は、論理/物理対応情報400を用いて論理領域に

対応する物理領域を求め、すなわち論理領域のアドレス（論理アドレス）を物理領域のアドレス（物理アドレス）に変換する（ステップ1010）。続いて制御部300は、リード/ライト対象は、この物理アドレスの記憶装置500からデータを読み出してホスト100に転送し、前記物理アドレスの記憶装置500に格納されたデータをリード/ライト対象は、ホスト100から転送されたデータを前記物理アドレスの記憶装置500に格納し（ステップ1020）、さらに後述の使用状況取得処理311を行う。リード/ライト要求およびデータ転送は1/0バス800を介して行われる。

【0029】論理/物理対応情報400の一例を図3に示す。論理アドレスはホスト100がリード/ライト処理310で用いる論理領域を示すアドレスである。物理アドレスは実際にデータが格納される記憶装置500上の領域を示すアドレスであり、記憶装置番号および記憶装置内アドレスからなる。記憶装置番号は個々の記憶装置500を示す。記憶装置内アドレスは記憶装置500内の記憶領域を示すアドレスである。

【0030】次に、使用状況取得処理311において制御部300は、リード/ライト処理310においてリード/ライト対象となった論理領域についての論理領域使用状況情報403と、リード/ライト処理310で使用する物理領域についての物理領域使用状況情報404を更新する（ステップ1030、1040）。論理領域使用状況情報403および物理領域使用状況情報404は、例えば使用頻度、使用率、リード/ライトに関する属性など、各々の論理領域と物理領域の各日時の使用状況に関する情報である。論理領域使用状況情報403および物理領域使用状況情報404の具体的な例は、以下の実施の形態で説明する。

【0031】次に、制御部300が行う再配置判断処理312について図4で説明する。

【0032】記憶装置500は、ユーザによって、また制御部300の分類（クラス600）に分類されており、クラス600への分類はクラス構成情報401によって設定されている。さらに、各クラス600は、ユーザによって、または初期条件として属性を設定されており、属性は、クラス属性情報402に設定されている。クラス属性情報402は、許容使用状況や好適な使用状況やクラス間優先順位などの属性に関する情報である。具体的な例は、以下の実施の形態で説明する。再配置判断対象期間情報405は、ユーザによってまたは初期条件として再配置判断処理312の対象とする使用状況情報と期間更新情報406が設定されている。

【0033】再配置判断対象期間情報405の一例を図5に示す。開始日時から終了日時までの期間が対象期間となる。期間更新情報は水回りの対象期間の設定条件であり、例えば毎週、毎日、X時間後などがあろう。制御部300は、対象期間の論理領域使用状況情報403および物理領域使用状況情報404を参照し、対象期間に

よび物理領域使用状況情報404を参照し（ステップ1100）、クラス属性情報402の各クラス600の許容使用状況などと比較して（ステップ1110）、物理的再配置を行うべき論理領域を選択する（ステップ1120）。

【0034】さらに、制御部300は、クラス属性情報402の許容使用状況や好適な使用状況やクラス間優先順位などを参照して（ステップ1130）、論理領域の再配置先（ステップ1140）を選択し（ステップ1150）、選択結果を再配置情報408に出力する（ステップ1160）。

【0035】再配置情報408の一例を図6に示す。論理領域は、再配置する論理領域であり、再配置先物理領域は、論理領域に対応する現在の物理領域を示す記憶装置番号と記憶装置内アドレスであり、再配置先物理領域は、再配置先の物理領域を示す記憶装置番号と記憶装置内アドレスである。図6に示すように再配置の立案は二つ以上行われる。さらに制御部300は、再配置判断対象期間情報405の期間更新情報を参照して、再配置判断対象期間情報405の対象期間を次回分を更新する（ステップ1170）。上記の処理において制御部300は、論理/物理対応情報400を用い、また前記の使用の物理領域の検索に未使用領域情報407を用いる。

【0036】未使用領域情報407の一例を図7に示す。記憶装置番号は個々の記憶装置500を示す。記憶装置内アドレスは記憶装置500内の領域を示すアドレスである。記憶装置番号および記憶装置内アドレスは物理領域を示し、使用/未使用の項目は、物理領域の使用/未使用の区別を示す。制御部300は、通常、再配置判断処理312を対象期間以後、後述の再配置実行処理313以前に自動的に行う。

【0037】次に、制御部300が行う再配置実行処理313について図8で説明する。

【0038】再配置実行時刻情報406はユーザによってまたは初期条件として再配置実行処理313を行う日時と日時更新情報406が設定されている。

【0039】再配置実行時刻情報406の一例を図9に示す。制御部300は、設定された日時以下に該当する再配置実行処理313を自動的に実行する。日時更新情報は水回りの再配置実行処理313を行う日時の設定条件であり、例えば毎週、毎日、X時間後などがあろう。制御部300は、再配置情報408に基づき再配置先物理領域に格納している内容を再配置先物理領域にコピーする（ステップ1200）。さらに、コピー完了して再配置先物理領域の内容が全て再配置先物理領域に反映された時点で、制御部300は、論理/物理対応情報400上の再配置を行う論理領域に対応する物理領域

を再配置元物理領域から再配置先物理領域に変更する（ステップ1210）。

【0040】さらに、制御部300は、未使用物理領域470上の再配置先物理領域を使用し、再配置元物理領域を未使用に変更する（ステップ1220）。さらに、制御部300は、再配置実行時刻情報406の日時更新情報を参照して、再配置実行時刻情報406の日時を次回分を更新する（ステップ1230）。

【0041】ユーザまたは保守員は、制御部300が上記の処理で用いている各情報を、制御部300からネットワーク900を介して、またはホスト100からネットワーク900または1/0バス800を介して設定および確認すること、特に、再配置情報408を確認および設定して再配置案を修正や追加や削除などをすることができ、

【0042】上記の処理を行うことにより、取得した使用状況情報および設定されたクラス属性に基づいて、ストレージサブシステム200において論理領域の物理的再配置を自動的に行い、ストレージサブシステム200の最適化を行うことができる。さらに上記の再配置判断および実行の処理を繰り返して配置を修正していくことによって、使用状況の変動やその他の最適化要因を吸収していくことができる。

【0043】特に、上記の処理により、ユーザまたは保守員は再配置による最適化を簡便に行うことができる。ユーザまたは保守員は、記憶装置500をクラス600という単位で管理できるため、記憶装置500の性能や信頼性や特性などの属性を個々の記憶装置500について管理する必要がある。さらに、ユーザまたは保守員は、必要に応じて同一の属性を持つクラス600に対して、1つの管理単位として扱うことができる。ただし、1つの記憶装置500が1つのクラス600を構成すると見なして1つの記憶装置500を管理単位として上記の再配置の処理を行うことも可能である。

【0044】また、ユーザまたは保守員は、ホスト100で行われる処理（ジョブ）の格好やスケジューリングに、計算機システムで行われる処理と、この処理に伴う1/0は、ユーザによって作成されたスケジューリングが行われる。ユーザは、特に最適化の対象とした処理を有する場合、処理の期間を特定することが可能であり、本実施の形態で説明した再配置の処理によって、ユーザは個々の物理領域を指定して再配置の処理を、ストレージサブシステム200に行わせ、すなわち、前記期間の使用状況情報に基づいて上記の再配置による最適化を実現することができる。また、計算機システムで行われる処理および1/0の傾向は日毎、月毎、年毎などの周期性を示す場合も多い。特に、処理が定常業務に基づく処理である場合には、周期性が顕著となる。前述の場

合と同様にユーザは、周期において特に最適化対象として個々の物理領域を指定して再配置による最適化を行うことができる。また、再配置実行処理313では、ストレージサブシステム200内で格納内容のコピーを行う。ユーザはストレージサブシステム200があまり使用されない時刻やホスト100で実行されている処理の要求処理性能が低い期間を再配置実行処理313の実行時刻として設定することで、ホスト100での要求処理性能が高い処理のストレージサブシステム200への1/0がコピーにより阻害されることを回避できる。

【0045】なお、記憶装置500は、それぞれ異なる性能、信頼性、特性や属性を持つていてよく、特に具体的には、磁気ディスク装置、磁気テープ装置、半導体メモリ（キャッシュ）のように異なる記憶媒体であってもよい。また、上記の例では未使用領域情報407は物理領域に基づいて記述されているとしたが、未使用の物理領域に対応する論理領域（論理アドレス）に基づいて記述されていてもよい。

【0046】＜第2の実施の形態＞本実施の形態では、使用状況情報としてのディスク装置使用率の適用と、クラス600の上限値およびクラス600間の性能順位による再配置判断について説明する。

【0047】図10は、本発明の第2の実施の形態における計算機システムの構成図である。

【0048】本実施の形態の計算機システムは、ホスト100、ディスク装置500、制御部300、制御部300を有している。本実施の形態における計算機システムは、第1の実施の形態でのストレージサブシステム200をディスク装置500と201とし、記憶装置500をパーティティグルーブ501としたものに相当する。

【0049】ディスク装置501は、制御部300とディスク装置502を有する。制御部300は、第1の実施の形態での制御部300に相当する。ディスク装置502は、n台（nは2以上の整数）でRAID（ディスクアレイ）を構成しており、このn台のディスク装置502による組をパーティティグルーブ501と呼ぶ。RAIDの性質として、1つのパーティティグルーブ501に含まれるn台のディスク装置502は、n-1台のディスク装置502の格納内容から生成される冗長データが残り1台に格納されるという冗長性上の関係を有する。またn台のディスク装置502は、冗長データを格納する。この関係から各パーティティグルーブ501を動作上の1単位とみなすことができるが、冗長性や台数などにより実現するためのコストや性能特性などが異なるため、ディスク装置502を構成するアレイ（パーティティグルーブ501）を混在させることも多く、またパーティティグルーブ501を構成するディスク装置502につ

いても、性能や容量などによりコストが異なるため、ディスクアレキシステム201を構成するにあたって最適なコストパフォーマンスを実現するために性能や容量の異なる複数のディスク装置502を用いることもあ
る。よって本実施の形態においてディスクアレキシステム201を構築する各パーティグループ501は性能、容量、特性などの属性が同一であるとは限らず、特に性能について差異があるとする。

【0050】本実施の形態における論理/物理対応情報400の一例を図11に示す。

【0051】論理アドレスは、ホスト100がリード/ライト処理310で用いる論理領域を示すアドレスである。物理アドレスは実際にデータと前記冗長データが格納されるディスク装置502上の領域を示すアドレスであり、パーティグループ番号と各々のディスク装置番号およびディスク装置内アドレスからなる。パーティグループ番号は個々のパーティグループ501を示す。ディスク装置番号は個々のディスク装置502を示す。ディスク装置内アドレスはディスク装置502内の領域を示すアドレスである。制御部300は、RAIDの動作として、冗長データに関する情報を前記リード/ライト処理310などで用いて処理するが、本実施の形態の図明では、パーティグループ502を動作上の1単位とし、説明する。前記処理に関してはここでは特に示さない。

【0052】さらに第1の実施の形態と同様に、パーティグループ501は、ユーザによってまたは初期状態として複数の組(クラス600)に分類されておられ、クラス600への分類はクラス構成情報401に設定されている。クラス構成情報401の一例を図12に示す。【0053】クラス番号は各クラス600を示す番号である。パーティグループ数は各クラス600に属するパーティグループの数を示す。パーティグループ番号は各クラス600に属するパーティグループ番号501を示す。同様に、各クラス600の属性は、クラス属性情報402に設定されている。本実施の形態におけるクラス属性情報402の一例を図13に示す。

【0054】クラス番号は、各クラス600を示す番号である。使用上限は後述のディスク使用率の許容範囲を示す上限であり、クラス600の属するパーティグループ501に適用する。クラス間性能順位は、クラス600間の性能順位(数字の小さいものが高性能とす)である。クラス間性能順位は各クラス600を構成するパーティグループ501の前述の性能属性に基づき、所配実行上限および固定については後述する。【0055】本実施の形態における使用状況情報取得処理311について図14で説明する。

【0056】制御部300は、第1の実施の形態と同様に、リード/ライト処理310において使用したディスク装置502の使用時間取得して単位時間当たりの使

用時間(使用率)を求め、さらに、ディスク装置502が属するパーティグループ501について、使用率の平均を算出し(ステップ1300)、使用率平均を、リード/ライト対象となった論理領域についてディスク装置使用率として論理領域使用状況情報403に記録する(ステップ1310)。また制御部300は、パーティグループ501に対応する全論理領域のディスク装置使用率の和を求め(ステップ1320)、パーティグループ501の使用率として物理領域使用状況情報404に記録する(ステップ1330)。

【0057】本実施の形態における論理領域使用状況情報403および物理領域使用状況情報404の一例を図15および図16に示す。

【0058】日時とはサンプリング間隔(一定期間)毎の日時を示し、論理アドレスは論理領域を示し、パーティグループ番号は個々のパーティグループを示し、論理領域のディスク装置使用率およびパーティグループ使用率はそれぞれ前記サンプリング間隔での平均使用率を示す。上記のようなディスク装置502の使用率は、ディスク装置502にかかると負荷を示す値であり、使用率が大きい場合は、ディスク装置502が性能ボトルネックとなっている可能性があるため、再配置処理で使用率を下げることによりディスクアレキシステム201の性能向上が期待できる。

【0059】次に、再配置判断処理312について図17で説明する。

【0060】制御部300は、各クラス600について、クラス600に属するパーティグループ501をクラス構成情報401から取得する(ステップ1300)。続いて、制御部300は、第1の実施の形態と同様の再配置判断対象期間情報405を参照して対象期間を取得し、さらにパーティグループ501について、対象期間の物理領域使用状況情報404のパーティグループ使用率を取得し集計する(ステップ1320)。続いて、制御部300は、クラス属性情報402を参照してクラス600の使用率上限値を取得する(ステップ1330)。制御部300は、パーティグループ使用率とクラス上限値と比較し、パーティグループ使用率がクラス上限値より大きい場合は、パーティグループ501の使用率を減らすために、パーティグループ501に対応する論理領域の再配置が必要と判断する(ステップ1340)。

【0061】続いて、制御部300は、対象期間の論理領域使用状況情報403を参照して、再配置が必要と判断したパーティグループ501の各物理領域に対応する論理領域のディスク装置使用率を取得し集計し(ステップ1350)、ディスク装置使用率の大きいものから、再配置する論理領域として選択する(ステップ1360)。論理領域の選択は、パーティグループ501の使用率から選択した論理領域のディスク装置使用率を減算し

ていき、クラス600の使用率上限値以下になるまで行う(1370)。ディスク装置使用率の大きい論理領域は、パーティグループ501の使用率に対する影響も大きく、またホスト100からの論理領域に対するアクセス頻度も大きいと考えられるため、ディスク装置使用率の大きい論理領域を優先的に再配置することで、ディスクアレキシステム201の効果的な性能改善が期待できる。

【0062】制御部300は、選択された論理領域についての物理領域となる物理領域を探し、パーティグループ501の属するクラス600より性能順位が高位のクラス600(高性能クラス)に注目し、クラス構成情報401および第1の実施の形態と同様の未使用物理領域情報407を参照して高性能クラスに属するパーティグループ501の未使用物理領域を取得する(ステップ1380)。

【0063】さらに、制御部300は、各未使用物理領域について、再配置先とした場合のパーティグループ使用率の平均値を求め(ステップ1390)、未使用物理領域の中から、再配置先とした場合に高性能クラスに設定されている上限値を超えないと予測できる未使用物理領域を、再配置先物理領域として選択し(ステップ1400)、選択結果を第1の実施の形態と同様に、再配置情報408に出力する(ステップ1410)。選択した全ての論理領域について再配置先の物理領域を選択し終えたら処理を終了する(ステップ1420)。

【0064】本実施の形態において、制御部300は、第1の実施の形態に加えてパーティグループ情報409を保持し、パーティグループ情報409、論理領域使用状況情報403、及び物理領域使用状況情報404から使用率予測値を算出する。

【0065】パーティグループ情報409の一例を図18に示す。パーティグループ番号は個々のパーティグループ501を示す番号である。RAID構成はパーティグループ501が構成するRAIDのレベルやディスク台数や冗長度構成を示す。ディスク装置性能はパーティグループ501を構成するディスク装置502の性能特性を示す。固定については後述する。上記の処理においてディスク装置使用率の大きい論理領域の再配置先を高性能クラスのパーティグループ501とする一方で、同一負荷に対するディスク装置使用率を短縮でき、論理領域の再配置後のディスク装置使用率を抑制できる。

【0066】再配置実行処理313は、第1の実施の形態と同様に行われるが、図19に示すように、制御部300は、再配置のためのコピーを行う前にクラス属性情報402を参照し、再配置先および再配置元のクラス600について、ユーザによってまたは初期条件として設定された再配置実行上限値を取得する(ステップ1500)。さらに物理領域使用状況情報404を参照して、

再配置先および再配置元のパーティグループ501の直近のパーティグループ使用率を取得し(ステップ1510)、比較の結果少なくとも一方のクラス600においてパーティグループ使用率が再配置実行上限値を超えていた場合は(ステップ1520、1530)、再配置実行処理313を中止または延期する(ステップ1540)。

【0067】上記処理によりユーザは、パーティグループ501の使用率が大きくならずながら負荷が高い場合に前記コピーによりさらに負荷が生じることを回避することができ、また回避のための上限値をクラス600毎に任意に設定することができる。

【0068】上記のように処理することによって、ディスク装置502の使用状況に基づいて物理的に再配置する論理領域の選択、および再配置先の物理領域の選択を、クラス構成および属性に基づいて行い、再配置によりディスク装置502の負荷を分散して、各クラス600に設定されている使用率上限値を、クラス600に属するパーティグループ501の使用率が超えない配置を実現することができ、さらに再配置判断および実行の処理を繰り返して配置を修正していくことによって、使用状況の変動や予測誤差を吸収していくことができる。

【0069】再配置判断処理312において、制御部300は、対象期間の物理領域使用状況情報404のパーティグループ使用率や、論理領域使用状況情報403の論理領域のディスク装置使用率を参照して集計し、判断に用いるとしたが、例えば、対象期間の全ての値の平均を用いる代わりに、対象期間中の上位m割の値を用いる方法も考えられ、また上位m番目の値を用いる方法も考えられる(mは1以上の整数)。これらの方法をユーザが選択できるようにすることで、ユーザは使用状況の特性的な部分のみを選択して用い、再配置判断処理312を行わせることができる。

【0070】上記の再配置判断処理312において、制御部300は、ディスクアレキシステム201の全てのクラス600について、論理領域の再配置の必要ないパーティグループ501の検出を行うとしたが、前記検出の前に制御部300がクラス属性情報402を参照し、固定属性が設定されているクラス600については、検出の対象外としてもよい。また同様に、制御部300がパーティグループ情報409を参照し、固定属性が設定されているパーティグループ501については検出の対象外としてもよい。また、再配置判断処理313において、制御部300は、高性能クラスに属するパーティグループ501の未使用物理領域から再配置先の物理領域を選択したとが、固定属性が設定されているクラス600については対象外として、さらに性能順位が高位のクラス600を高性能クラスとして、さらに性能順位が高位のクラス600を高性能クラスとして扱うようにしてもよい。また固定属性が設定されているパーティグループ501については対象外としてもよい。上記のように図

定属性が設定されているクラス600またはパリティグループ501を扱うことにより、ユーザは上記の自動的な再配置処理において物理的な再配置の影響を生じさせなく、クラス600またはパリティグループ501を指定し、再配置の対象とすることができ、

【0071】＜第3の実施の形態＞本実施の形態では、同一クラス600内での再配置判断について説明する。本実施の形態での計算機システムは、第2の実施の形態と同様である。ただし、本実施の形態では1つのクラス600に複数のパリティグループ501が属する。本実施の形態での処理は、再配置判断処理312を除いては第2の実施の形態と同様である。また、再配置判断処理312についても、再配置する論理領域の選択（ステップ1600）は、第2の実施の形態と同様である。

【0072】本実施の形態での再配置判断処理312における、再配置先の物理領域の選択について図20で説明する。

【0073】第2の実施の形態では再配置先の物理領域を再配置元の物理領域の属するクラス600より性能順位が低いクラスのクラス600から選択するが、本実施の形態では同一クラス600の再配置元以外のパリティグループ501から選択する。制御部300は、クラス構成情報401と未使用領域情報407を参照して、同一クラス600に属する再配置元以外のパリティグループ501の未使用物理領域を取得する（ステップ1610）。

制御部300は、各未使用物理領域について、再配置先とした場合のパリティグループ501の使用率の予測値を求め（ステップ1620）、未使用物理領域の中から、再配置先とした場合に同一クラス600に設定されている上乗積を超えないと予測できる未使用物理領域を、再配置先の物理領域として選択し（ステップ1630）、選択結果を第2の実施の形態と同様に、再配置情報408に出力する（ステップ1640）。再配置する全ての論理領域について再配置先の物理領域を選択し終えたら処理を終了する（ステップ1650）。

【0074】上記の処理により、同一クラス600内においてディスク装置502の負荷を分散することができ、上記の処理方法は例えばディスクアレイシステム201のパリティグループ501が全て1つのクラス600（単一クラス）に属する構成に適用することができ、また、例えば、第2の実施の形態で説明した処理方法と組み合わせた場合に、再配置先の未使用物理領域の選択において、再配置元のクラス600より性能順位が低いクラス600に最適な未使用物理領域が得られない場合や、性能順位が最上位のクラス600での処理に適用できる。第2の実施の形態で説明した処理方法と組み合わせた場合は、第2の実施の形態での処理方法と異なる使用率上限値を用いてもよく、すなわち、そのためにクラス属性情報402が各クラス600について

二種類の使用率上限値または差分を有してもよい。

【0075】＜第4の実施の形態＞本実施の形態では、第2の実施の形態での再配置判断処理312において、再配置元のクラス600より性能順位が低いクラスのクラス600（高性能クラス）に再配置先の未使用物理領域が見つからなかつた場合に、再配置先を得るために先立って行われる、性能順位がより低いクラスのクラス600（低性能クラス）への高性能クラスからの再配置の処理について説明する。

【0076】本実施の形態での計算機システムは、第2の実施の形態と同様である。本実施の形態における再配置判断処理312について図21で説明する。

【0077】制御部300は、高性能クラスに属するパリティグループ501をクラス構成情報401から取得する（ステップ1700）。続いて制御部300は、第1の実施の形態と同様の再配置判断対象期間情報405を参照して対象期間を取得し（ステップ1710）、対象期間の論理領域の使用状況情報403を参照して、パリティグループ501の未使用物理領域に属する論理領域のディスク装置使用率を取得し（ステップ1720）、ディスク装置使用率の小さいものから、低性能クラスへ再配置する論理領域として選択する（ステップ1730）。このとき論理領域の選択は必要だけ行われる（ステップ1740）。

【0078】続いて制御部300は、選択された論理領域についての再配置先となる物理領域を、低性能クラスに属するパリティグループ501から選択するが、再配置先の物理領域選択の処理は、第2の実施の形態での処理説明において再配置先としている高性能クラスを低性能クラスと置き換えれば、第2の実施の形態での処理と同様である（ステップ1750）。また、本実施の形態におけるその他の処理も第2の実施の形態での処理と同様である。

【0079】上記の処理を行うことで、第2の実施の形態での再配置判断処理312において高性能クラスに再配置先の未使用物理領域が見つからなかつた場合に、高性能クラスから低性能クラスへ論理領域の再配置を、高性能クラスへの再配置に先立って行い、再配置先の未使用物理領域を高性能クラスに用意することができ、制御部300は、上記の処理を必要に応じて繰り返して、十分な未使用物理領域を用意することができ、

【0080】論理領域の再配置先を低性能クラスのパリティグループ501とするため、同一負荷に対するディスク使用時間が再配置について増大し、論理領域の再配置後のディスク装置使用率が增大する可能性があるが、ディスク使用率の小さい論理領域から再配置していくことで、増大の影響を最小限に抑えることができる。

【0081】＜第5の実施の形態＞本実施の形態では、クラス600の属性の1つにアクセス種別属性を設け、

アクセス種別属性を用いてシーケンシャルアクセスが頻りに行われる論理領域とランダムアクセスが頻りに行われる論理領域とを、他のパリティグループ501に自動的に物理的に再配置して分離するための再配置判断について説明する。

【0082】本実施の形態における計算機システムは図10に示したものである。本実施の形態では、第2の実施の形態での説明に加え、制御部300が保持する下記の情報を用いる。

【0083】本実施の形態でのクラス属性情報402の一例を図22に示す。この例では、第2の実施の形態の例に対しアクセス種別が加えられており、クラス600のアクセス種別が、例えばシーケンシャルに設定されている場合は、クラス600がシーケンシャルアクセスに好適であると設定されていることを示す。

【0084】本実施の形態での論理領域使用状況情報403の一例を図23に示す。この例では、第2の実施の形態の例に対し、シーケンシャルアクセス事およびランダムアクセス事が加えられている。

【0085】さらに、本実施の形態において制御部300は、第2の実施の形態に加え、アクセス種別属性情報410と論理領域属性情報411を保持する。

【0086】アクセス種別属性情報410の一例を図24に示す。ユーザによりまたは初期条件として、アクセス種別属性情報411の一例を図25に示す。アクセス種別属性情報は、各論理領域について頻りに行われると期待できるアクセス種別であり、ユーザが設定する。固定については後述する。

【0087】本実施の形態での処理は、使用状況情報取得処理311および再配置判断処理312を除いては第2の実施の形態と同様である。

【0088】本実施の形態における使用状況情報取得処理311について図26で説明する。

【0089】制御部300は、第2の実施の形態での使用状況情報取得処理311と同様に、論理領域についてのディスク装置使用率を算出し（ステップ1800、1810）、リード/ライト処理310の使用率内容を分析して、使用率についてシーケンシャルアクセスとランダムアクセスの比率を算出し（ステップ1820）、使用率およびアクセス種別比率を論理領域使用状況情報403に記録する（ステップ1830）。また、制御部300は、第2の実施の形態と同様にパリティグループ501の再配置先と物理領域使用状況情報404への記録を行う（ステップ1840、1850）。

【0090】本実施の形態における再配置判断処理312において、再配置する論理領域の選択は第2の実施の形態と同様である（ステップ1900）。再配置判断処理312での再配置先の物理領域の選択について図27

で説明する。

【0091】制御部300は、論理領域使用情報403を参照し、再配置する論理領域についてのシーケンシャルアクセス率を取得し（ステップ1910）、アクセス種別属性情報410に設定されている基準値と比較する（ステップ1920）。シーケンシャルアクセス率が基準値より大きい場合、制御部300は、クラス属性情報402を参照し、アクセス種別がシーケンシャルと設定されているクラス600（シーケンシャルクラス）と存在する（ステップ1950）。シーケンシャルクラスが存在する場合、制御部300は、クラス構成情報401と未使用物理領域情報407を参照して、シーケンシャルクラスに属する再配置元以外のパリティグループ501の未使用物理領域を取得する（ステップ1960）。さらに制御部300は、各未使用物理領域について、再配置先とした場合のパリティグループ501の使用率の予測値を求め（ステップ1970）、未使用物理領域の中から、再配置先とした場合にシーケンシャルクラスに属する論理領域を超えないと予測できる未使用物理領域を、再配置先の物理領域として選択し（ステップ1980）、選択結果を第2の実施の形態と同様に再配置情報408に出力する（ステップ1990）。制御部300は、使用率予測値を、第2の実施の形態と同様のパリティグループ情報409と本実施の形態における論理領域使用状況情報403および物理領域使用状況情報404から算出する。

【0092】上記の比較において、シーケンシャルアクセス率は基準値以下である場合、制御部300は、論理領域属性情報411を参照し、論理領域についてアクセス種別属性がシーケンシャルと設定されているか調べ（ステップ1940）。アクセス種別属性にシークンシャルと設定されている場合、上記と同様に制御部300は、シーケンシャルクラスの有無を調べ（ステップ1950）、シーケンシャルクラスが存在する場合は、シーケンシャルクラスから再配置先の物理領域を選択する（ステップ1960～1990）。

【0093】上記の比較において、シーケンシャルアクセス率が基準値以下であり、さらにアクセス種別属性がシーケンシャルでなかつた場合、またはシーケンシャルクラスが存在しなかつた場合、制御部300は、第2の実施の形態と同様に、シーケンシャルクラスの物理領域を選択する（ステップ2000）。

【0094】上記の処理により、同一パリティグループ501での異なるシーケンシャルアクセスとランダムアクセスの混在に対し、各クラス600に属性として設定されたアクセス種別と使用率上限値を用いて、ランダムアクセスとシークンシャルとを行われる論理領域とを、異なるパリティグループ501に自動的に再配置して分離、すなわち異

なるディスク装置502に分離することができ、特にランダムアクセスに対する応答性を改善することができる。

【0095】また、上記の処理においては制御部300は、シーケンシャルアクセスに注目して再配置による自動的分離を行うとしたが、同様ランダムアクセスに注目して分離を行うことも可能である。

【0096】上記の再配置処理312において、再配置する論理領域を選択した時点で、制御部300が論理領域属性情報411を参照し、論理領域に固有属性が指定されている場合は、論理領域を再配置しないとするが、ユーザが特に再配置を行いたくないと考える論理領域がある場合、固有属性を設定することで論理領域を再配置の対象とすることができ、上記の固有属性に関する処理は論理領域属性情報411を用いることで、前述の実施の形態にも適用できる。

【0097】

【発明の効果】ストレージサブシステムのユーザ、または保守員が、記憶領域の物理的再配置による配置最適化を行うための作業を簡便にすることができ、

【図面の簡単な説明】

【図1】本発明の第1の実施の形態での計算機システムの構成図である。

【図2】本発明の第1の実施の形態でのリード/ライト処理310および使用状況取得処理311のフローチャートである。

【図3】本発明の第1の実施の形態での論理/物理対応情報400の一例を示す図である。

【図4】本発明の第1の実施の形態での再配置判断処理312のフローチャートである。

【図5】本発明の第1の実施の形態での再配置判断対象期間情報405の一例を示す図である。

【図6】本発明の第1の実施の形態での再配置情報408の一例を示す図である。

【図7】本発明の第1の実施の形態での未使用領域情報407の一例を示す図である。

【図8】本発明の第1の実施の形態での再配置実行処理313のフローチャートである。

【図9】本発明の第1の実施の形態での再配置実行時刻情報406の一例を示す図である。

【図10】本発明の第2の実施の形態および第5の実施の形態の計算機システムの構成図である。

【図11】本発明の第2の実施の形態での論理/物理対応情報400の一例を示す図である。

【図12】本発明の第2の実施の形態でのクラス構成情報401の一例を示す図である。

【図13】本発明の第2の実施の形態でのクラス属性情報402の一例を示す図である。

【図14】本発明の第2の実施の形態での使用状況取得処理311のフローチャートである。

【図15】本発明の第2の実施の形態での論理領域使用状況情報403の一例を示す図である。

【図16】本発明の第2の実施の形態での物理領域使用状況情報404の一例を示す図である。

【図17】本発明の第2の実施の形態での再配置判断処理312のフローチャートである。

【図18】本発明の第2の実施の形態でのパリティグループ情報409の一例を示す図である。

【図19】本発明の第2の実施の形態での再配置実行処理313のフローチャートである。

【図20】本発明の第3の実施の形態での再配置判断処理312のフローチャートである。

【図21】本発明の第4の実施の形態での再配置判断処理312のフローチャートである。

【図22】本発明の第5の実施の形態でのクラス属性情報402の一例を示す図である。

【図23】本発明の第5の実施の形態での論理領域使用状況情報403の一例を示す図である。

【図24】本発明の第5の実施の形態でのアクセス種別基礎値情報410の一例を示す図である。

【図25】本発明の第5の実施の形態での論理領域属性情報411の一例を示す図である。

【図26】本発明の第5の実施の形態での使用状況取得処理311のフローチャートである。

【図27】本発明の第5の実施の形態での再配置判断処理312のフローチャートである。

【符号の説明】

100 ホスト

200 ストレージサブシステム

201 ディスクアレイシステム

300 制御部

310 リード/ライト処理

311 使用状況取得処理

312 再配置判断処理

313 再配置実行処理

400 論理/物理対応情報

401 クラス構成情報

402 クラス属性情報

403 論理領域使用状況情報

404 物理領域使用状況情報

405 再配置判断対象期間情報

406 再配置実行時刻情報

407 未使用領域情報

408 再配置情報

409 パリティグループ情報

410 アクセス種別基礎値情報

411 論理領域属性情報

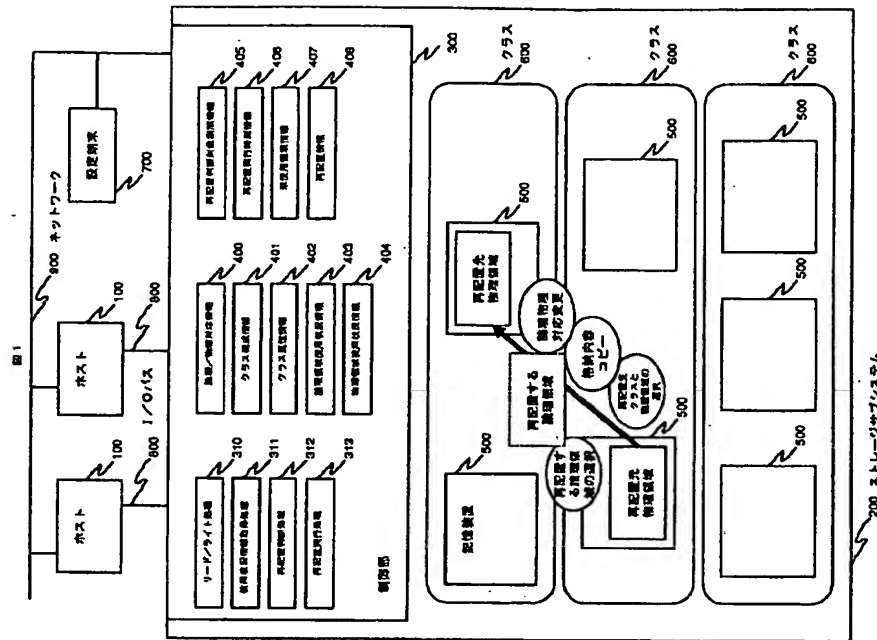
500 記憶装置

501 パリティグループ

502 ディスク装置

600 クラス
700 制御端末
800 I/Oバス
900 ネットワーク

【図1】



【図9】

【図24】

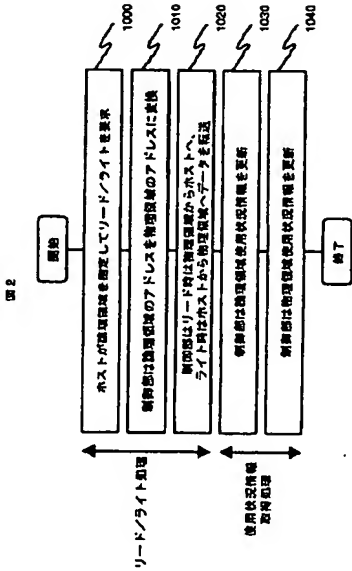
図9

図24

日時	1999年8月11日 22時0分
日時更新情報	毎日 (+24時間)

7フセム型基礎値 (%)	75
--------------	----

【図2】



【図3】

図3

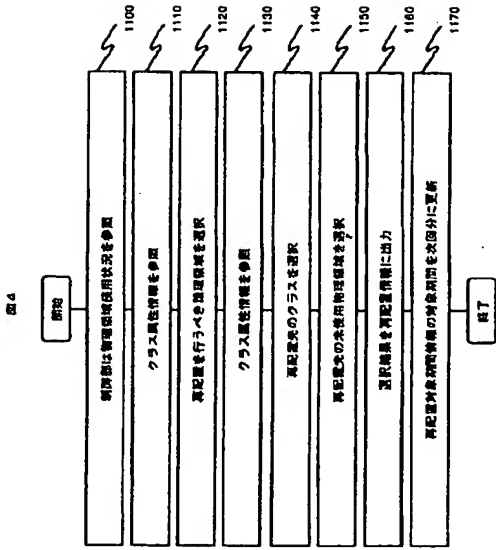
物理アドレス	物理アドレス	
	記憶装置番号	記憶装置内アドレス
0~999	0	0~999
1000~1999	0	1000~1999
2000~2999	1	0~999
3000~3999	1	1000~1999

【図5】

図5

開始日時	1999年8月11日 0時30分
終了日時	1999年8月11日 17時15分
時間更新情報	毎日 (+24時間)

【図4】



【図6】

図6

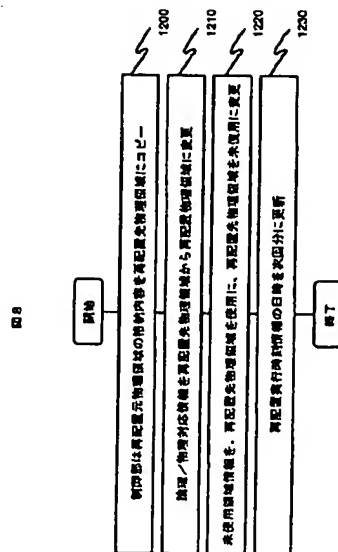
番号	記憶装置	物理部属性情報		物理部属性情報	
		記憶装置番号	記憶装置内アドレス	記憶装置番号	記憶装置内アドレス
1	0~999	0	0~999	10	0~999
2	1000~1999	0	1000~1999	10	1000~1999

【図7】

図7

記憶装置番号	記憶装置内アドレス	使用/未使用
0	0~999	使用
0	1000~1999	使用
0	2000~2999	未使用
0	3000~3999	未使用

【图8】



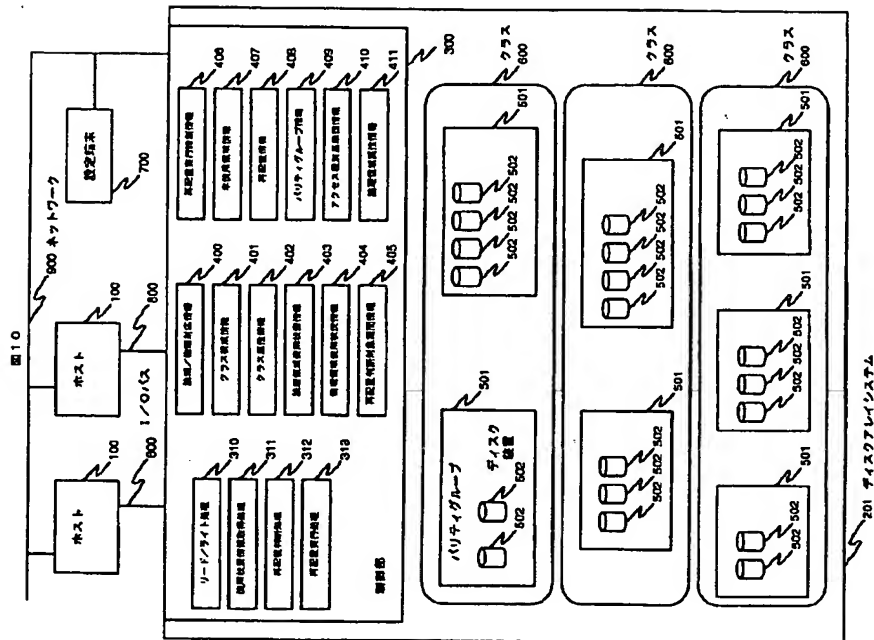
【圖 11】

姓 名 アドレス	所 属 アドレス				
	パリティグループ 番号	データ		冗長データ	
		記憶装置 番号	記憶装置内 アドレス		
0～999	100	0	0～999	20	0～999
1000～1999	100	0	1000～1999	20	1000～1999
2000～2999	101	1	0～999	41	0～999
3000～3999	101	1	1000～1999	41	1000～1999

【图 12】

クラス番号	パリティグループ数	パリティグループ番号
0	3	100, 110, 120
1	2	101, 111
2	4	102, 112, 122, 132

【圖10】



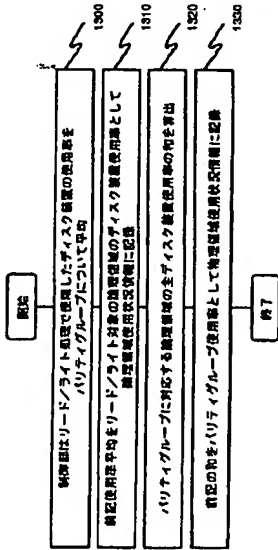
【図13】

図13

クラス番号	使用率上限値 (%)	クラス属性記号	対応実行上限値 (%)	固定
0	60	1	70	-
1	70	2	80	固定
2	80	3	90	-

【図14】

図14



【図15】

図15

日時	クラス属性記号	ディスク使用率 (%)
1999年8月11日 8時0分	0-999	18
1999年8月11日 8時15分	1000-1999	32
1999年8月11日 8時30分	0-999	20
1999年8月11日 8時45分	1000-1999	30
1999年8月11日 9時0分	0-999	22
1999年8月11日 9時15分	1000-1999	28

【図16】

図16

日時	バリエーション番号	使用率 (%)
1999年8月11日 8時0分	100	68
1999年8月11日 8時15分	101	52
1999年8月11日 8時30分	100	70
1999年8月11日 8時45分	101	50
1999年8月11日 9時0分	100	72
1999年8月11日 9時15分	101	48

【図18】

図18

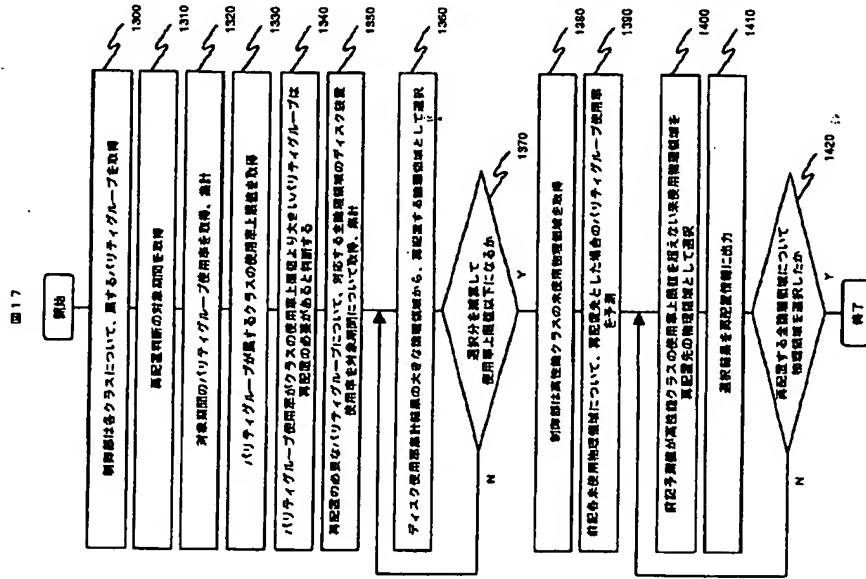
バリエーション番号	RAID構成	ディスク使用率 (%)	固定
100	RAID0/1	110	-
101	RAID1/0	100	固定
102	RAID0/1	95	-

【図22】

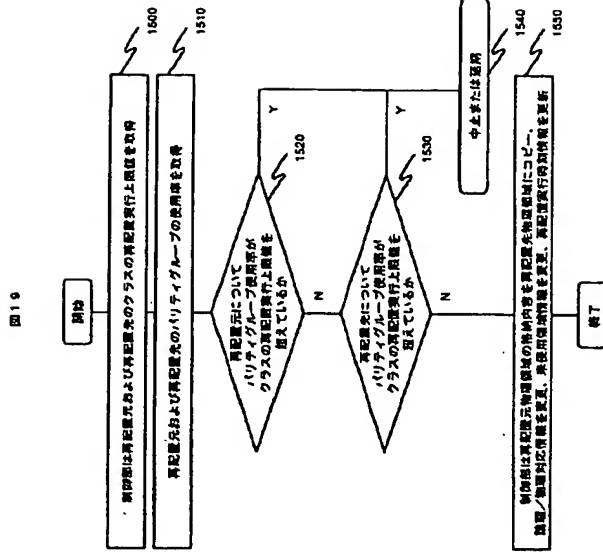
図22

クラス番号	使用率上限値 (%)	クラス属性記号	対応実行上限値 (%)	固定	アクセス場所
0	60	1	70	-	-
1	70	2	80	-	-
2	80	3	90	-	ランダムアクセス

【図17】



【図19】

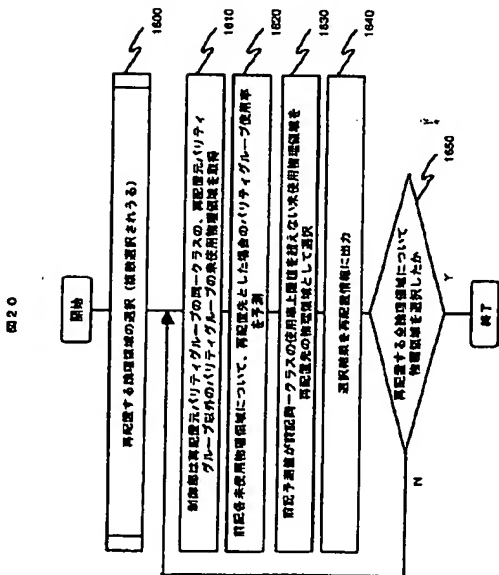


【図23】

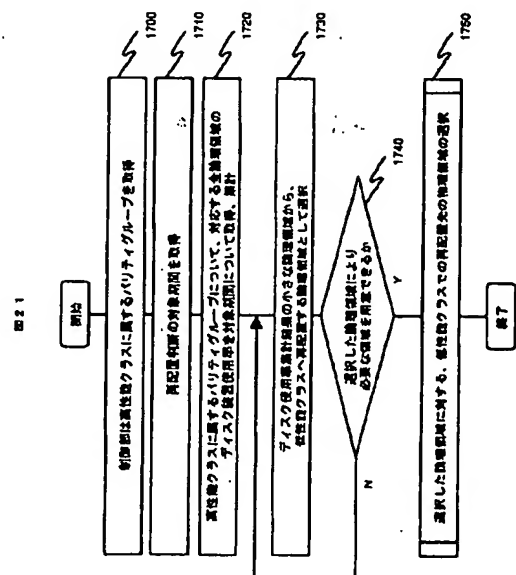
図23

日時	始期アドレス	ディスク容量 使用率 (%)	シーケンシャル アクセス率 (%)	ランダム アクセス率 (%)
1999年8月11日 8時0分	0~999	18	78	22
	1000~1999	32	82	49
1999年8月11日 8時15分	0~999	20	80	20
	1000~1999	30	80	50
1999年8月11日 8時30分	0~999	22	82	18
	1000~1999	28	48	52

【図20】



【図21】

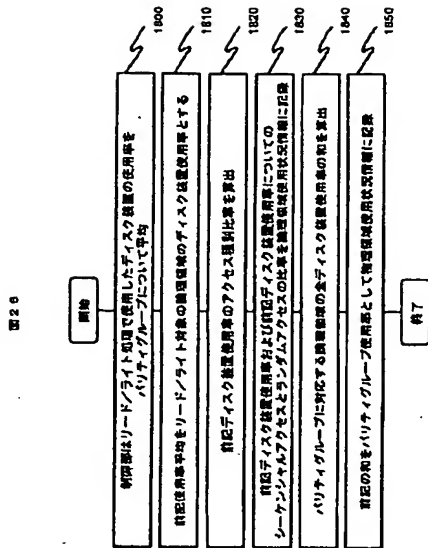


【図25】

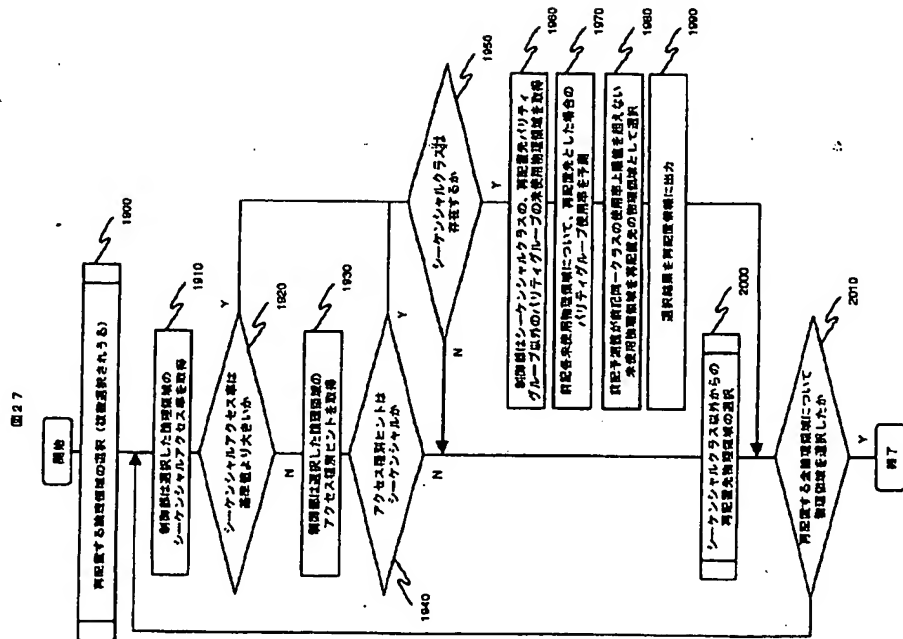
図25

管理アドレス	アクセス頻度	管理
0~999	-	-
1000~1999	-	-
2000~2999	シーケンシャル	-
3000~3999	-	固定

【図26】



【図27】



フロントページの続き

(72)発明者 山神 忍司
神奈川県川崎市麻生区王禅寺109番地 株
式会社日立製作所システム開発研究所内

(72)発明者 荒井 弘治
神奈川県小田原市国府津2880番地 株式会
社日立製作所ストレージシステム事業部内